

. Alonso A, Seguí-Gómez M, De Irala J, Sánchez Villegas A, Beunza JJ, Martínez-González MA. Predictors of follow-up and assessment of selection bias from dropouts using inverse probability weighting in a cohort of university graduates. *Eur J Epidemiol* 2006;21(5):351-358. PMID 16736275.

2 **Predictors of follow-up and assessment of selection bias from dropouts**
3 **using inverse probability weighting in a cohort of university graduates**

5 Alvaro Alonso^{1,2}, María Seguí-Gómez^{1,3}, Jokin de Irala¹, Almudena Sánchez-Villegas^{1,4},
6 Juan José Beunza¹ & Miguel Ángel Martínez-Gonzalez¹

7 ¹Department of Preventive Medicine and Public Health, Clínica Universitaria, University of Navarra, Pamplona, Spain;
8 ²Department of Epidemiology, Harvard School of Public Health, Boston, MA, USA; ³Bloomberg School of Public Health,
9 Johns Hopkins University, Baltimore, MD, USA; ⁴Department of Clinical Sciences, University of Las Palmas de Gran
10 Canaria, Las Palmas de Canaria, Spain

11 Accepted in revised form 13 March 2006

12 **Abstract.** Dropouts in cohort studies can introduce
13 selection bias. In this paper, we aimed (i) to assess
14 predictors of retention in a cohort study (the SUN
15 Project) where participants are followed-up through
16 biennial mailed questionnaires, and (ii) to evaluate
17 whether differential follow-up introduced selection
18 bias in rate ratio (RR) estimates. The SUN Study
19 recruited 9907 participants from December 1999 to
20 January 2002. Among them, 8647 (87%) participants
21 answered the 2-year follow-up questionnaire. The
22 presence of missing information in key variables at
23 baseline, being younger, smoker, a marital status dif-
24 ferent of married, being obese/overweight and a his-
25 tory of motor vehicle injury were associated with being
26 lost to follow-up, while a self-reported history of car-
27 diovascular disease predicted a higher retention pro-
28 portion. To assess whether differential follow-up

affected RR estimates, we studied the association be- 29
tween body mass index and the risk of hypertension, 30
using inverse probability weighting (IPW) to adjust 31
for confounding and selection bias. Obese individuals 32
had a higher crude rate of hypertension compared 33
with normoweight participants (RR = 6.4, 95% 34
confidence interval (CI): 3.9–10.5). Adjustment for 35
confounding using IPW attenuated the risk of hyper- 36
tension associated to obesity (RR = 2.4, 95% CI: 37
1.1–5.3). Additional adjustment for selection bias did 38
not modify the estimations. In conclusion, we show 39
that the follow-up through mailed questionnaires of a 40
geographically disperse cohort in Spain is possible. 41
Furthermore, we show that despite existing differences 42
between retained or lost to follow-up participants this 43
may not necessarily have an important impact on the 44
RR estimates of HTN associated to obesity. 45

46 **Key words:** Attrition, Body mass index, Cohort studies, Hypertension, Inverse probability weighting, Selection
47 bias

48 **Abbreviations** BMI = body mass index; CI = confidence interval; HTN = hypertension; IPW = inverse
49 probability weighting; MET = metabolic equivalent; RR = rate ratio; SUN = Seguimiento Universidad de
50 Navarra

51 **Introduction**

52 Selection bias has been considered an especially
53 serious drawback in case-control studies, but this
54 problem can also afflict cohort studies. In the latter
55 study design, selection bias may arise from three
56 mechanisms: (i) unwillingness to participate, (ii)
57 missing information in some covariates (and thus,
58 exclusion from some analyses) and (iii) attrition of
59 the cohort (dropouts or losses to follow-up) [1, 2].
60 Dropouts can lead to selection bias because individ-
61 uals who do not respond in the follow-up tend to be
62 different from respondents [3]. This can be the case
63 particularly if the exposure of interest is associated
64 with censoring and there are variables associated
65 both with the study outcome and with censoring [4].

A first step in the study of a possible selection bias 66
introduced by cohort attrition is to compare the 67
characteristics of respondents and non-respondents. 68
Afterwards, different tools can be used to explore the 69
magnitude of bias and to correct it [5]. A recently 70
proposed approach to adjust for selection bias is 71
inverse probability weighting (IPW) [4, 6]. This meth- 72
od works by creating a pseudopopulation, in which 73
uncensored individuals are weighted by the inverse of 74
the probability of being non-censored given past 75
exposure and past levels of potential confounding 76
variables. In this way, the study participant accounts in 77
the analysis for those with similar characteristics that 78
were not selected (because they were censored). 79

Cohort studies can have a variety of designs. Some 80
are based in geographically dispersed populations 81

82 with special characteristics and their participants are
 83 usually followed-up through mailed questionnaires.
 84 This is an efficient way to carry out large prospective
 85 studies. However, most of these studies are being
 86 conducted in the United States [7-9], and there is no
 87 information about the follow-up of cohorts with
 88 these characteristics in Europe. Also, given the diffi-
 89 culty to track cohort participants exclusively through
 90 mail [10], it is valuable to assess factors associated
 91 with dropping out from the study.

92 The *Seguimiento Universidad de Navarra* (SUN,
 93 University of Navarra Follow-Up) Study was estab-
 94 lished in 2000 in Spain [11]. Participants in this cohort
 95 are all university graduates, recruited and followed-up
 96 through biennial mailed questionnaires. The objective
 97 of this paper was twofold: (i) to describe the initial
 98 follow-up of the SUN cohort and (ii) to investigate
 99 whether losses to follow up introduced bias in our
 100 estimations. For this latter objective, we studied the
 101 potential bias due to differential follow-up with a
 102 particular example, the association between body
 103 mass index (BMI) at baseline and the risk of incident
 104 hypertension (HTN) using IPW.

105 **Methods**

106 *The SUN Study*

107 The SUN Study is a multipurpose, dynamic cohort of
 108 university graduates in Spain. The SUN Study was
 109 approved by the Institutional Review Board of the
 110 University of Navarra. Among its objectives are the
 111 study of the association between dietary and other
 112 lifestyle variables and the incidence of cardiovascular
 113 disease, hypertension, obesity, and diabetes, and the
 114 assessment of risk factors for motor-related injuries.
 115 Its methods have been described previously [11].
 116 Briefly, beginning on December 1999, all university
 117 graduates from the University of Navarra, and
 118 university graduates from some professional associ-
 119 ations received a letter of invitation to participate in
 120 the study, a questionnaire to respond and a postage-
 121 prepaid envelope to return the questionnaire. This
 122 baseline questionnaire gathered information about
 123 sociodemographic variables, lifestyle factors, includ-
 124 ing physical activity, clinical variables and included a
 125 detailed food frequency questionnaire previously
 126 validated in Spain [12]. To quantify the volume and
 127 intensity of leisure-time physical activity, we calcu-
 128 lated an activity metabolic equivalent (MET) index
 129 for each participant. We assessed each participant's
 130 involvement and time spent in 17 different activities.
 131 We assigned a multiple of the resting metabolic rate
 132 (MET score) to each of these activities using previ-
 133 ously published guidelines [13]. The MET score of
 134 each activity was multiplied by the weekly time
 135 spent in each activity and a value of overall weekly
 136 MET-hours was obtained.

Follow-up of the cohort

In the baseline questionnaire, we requested three
 postal addresses from each participant (main address,
 alternative personal address, and name and address
 of a relative or friend) and at least one e-mail address.

The follow-up is conducted through biennial
 mailed questionnaires. Once a participant reaches its
 second (or fourth, or sixth) anniversary from the
 reception of the baseline questionnaire, we send him/
 her a follow-up questionnaire. To ensure high reten-
 tion rates, we send up to five reminders to non-
 respondents, the last of which is mailed-certified. In
 the non-questionnaire year (i.e., first, third, etc.), a
 salutatory letter is sent to all participants reminding
 them of the importance of continuous collaboration
 and the need to update their addresses.

Participants can contact the study investigators
 through email, telephone and postal address. Inter-
 estingly, e-mail has been an important and inexpen-
 sive mean to reach some participants and to
 communicate with them. Internet-based tools have
 been introduced to both update addresses and answer
 questionnaires.

In this paper we included all participants recruited
 up to January 2002 who answered the 2-year follow-up
 questionnaire before 25 July 2004.

Evaluation of bias: association between BMI and incident HTN

To assess the impact of selection bias introduced by
 losses to follow-up, we studied the association be-
 tween BMI in the baseline questionnaire, a classical
 cardiovascular risk factor, and the risk of incident
 HTN as assessed in the 2-year follow-up question-
 naire. We decided to study this association because,
 first, BMI was associated with the probability of
 retention in our study population and, two, because it
 is a well-known risk factor for HTN. BMI was
 computed as the ratio between baseline self-reported
 weight in kilograms divided by squared self-reported
 height (in meters). The repeatability and validity of
 self-reported weight have been the objective of a
 previous publication [14]. There, we showed that the
 mean relative error in self-reported weight was 1%.
 Also, we assessed the validity of the self-reported
 diagnosis of HTN in a random sample of 79 indi-
 viduals reporting a medical diagnosis of HTN and 41
 not reporting such diagnosis living in the metropoli-
 tan area of Pamplona (Navarra, Spain). These par-
 ticipants were representative of the SUN study
 population. Two study physicians performed two
 blood pressure measurements and an interview with
 these participants. The mean of both blood pressure
 measurements was computed to define hypertensive
 status. We defined HTN as systolic blood pressure
 ≥ 140 , or diastolic blood pressure ≥ 90 , or taking anti-
 HTN medication. Using this definition, the positive

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

159

160

161

162

163

164

165

166

167

168

169

170

171

172

173

174

175

176

177

178

179

180

181

182

183

184

185

186

187

188

189

190

191

192

193 and negative predictive value for the self-reported
194 diagnosis of HTN were 82 and 85%, respectively. A
195 detailed description of this validation study has been
196 published elsewhere [15].

197 For this particular evaluation we excluded other-
198 wise eligible individuals who reported HTN, cardio-
199 vascular disease, cancer or diabetes at baseline, those
200 with missing values in any covariate of interest, and
201 those with extremely low or high caloric intakes at
202 baseline (lower than 400 or higher than 3500 kcal/day
203 for women, and lower than 600 or higher than
204 4200 kcal/day for men). We consider as potential
205 confounding factors those variables previously
206 reported to be determinants of the incidence of HTN
207 (age, sex, physical activity, alcohol intake, fruit and
208 vegetable consumption), markers of a healthier life-
209 style (smoking), and variables we found associated
210 with dropout probability (marital status).

211 *Statistical analysis*

212 To attain the first objective, description of the follow-
213 up of the cohort, we compared respondents and non-
214 respondents to the first follow-up questionnaire using
215 logistic regression. The outcome was non-response to
216 the follow-up questionnaire, and as independent
217 candidate predictors we studied variables previously
218 associated to dropout in other studies, such as dif-
219 ferent sociodemographic variables, and risk factors
220 for cardiovascular disease, HTN, and motor vehicle
221 injuries (the main outcomes of interest in the SUN
222 study). Continuous variables were categorized for all
223 the analyses. Initially we ran univariate logistic
224 regression models for each independent variable.
225 Then, we ran a multivariate model simultaneously
226 including every variable considered in the univariate
227 analysis. In this analysis, we excluded individuals
228 with missing values in any variable to allow the model
229 to run properly. Goodness-of-fit was assessed with
230 the Hosmer-Lemeshow test.

231 For the second aim of the study, we used a marginal
232 structural Cox proportional hazards models to study
233 the association between BMI and incidence of HTN
234 [16], using IPW to adjust for confounding and selection
235 bias related to losses to follow-up. In these models, we
236 categorized the BMI as lower than 25 kg/m², from 25
237 through 30 kg/m², and greater than 30 kg/m² [17].

238 The use of IPWs in a marginal structural
239 model allows the estimation of causal effects in the

of conditional exchangeability in observational 243
studies see, for example, [18] and [6]. Briefly, we say 244
that an exposure E has a causal effect in a population 245
if the proportion of subjects developing a disease D 246
had everybody in the population been exposed to E is 247
different to the proportion of subjects developing that 248
disease had everybody been unexposed. In other 249
words, the exposure E has a causal effect if there is a 250
difference between the counterfactual risk of the 251
outcome D had everybody been exposed and the 252
counterfactual risk everybody had not been exposed. 253
By conditional exchangeability we mean that, after 254
conditioning on potential confounders, we can as- 255
sume that the probability of having a disease D had 256
everybody been exposed to E or had they been 257
unexposed is independent of the actual exposure 258
status, i.e. in every level of potential confounders the 259
exposure E can be assumed to be assigned randomly. 260
That is, conditional exchangeability implies that ex- 261
posed and unexposed individuals are exchangeable 262
across levels of potential confounders. 263

The IPW method to adjust for confounding works 264
by creating a pseudopopulation where each individ- 265
ual in the population is weighted by the inverse of the 266
conditional probability of receiving the exposure that 267
she or he received (conditional to the actual values of 268
the confounding variables). In that pseudopopulation 269
there is no confounding for the variable used to 270
condition for when calculating the weights, under the 271
assumption of no residual confounding. 272

To adjust for confounding, we computed stabilized 273
weights for each individual included in the final 274
analysis as the ratio of the probability of being in the 275
corresponding BMI category over the probability of 276
being in the observed BMI category for an individual 277
given her or his distribution of potential confounders 278
(Equation 1). Stabilized weights tend to be much less 279
variable than standard weights (as defined in the 280
previous paragraph) [16]. The numerator probability 281
could be computed directly from the data, as the 282
proportion of individuals in each BMI category. The 283
denominator was estimated using ordinal logistic 284
regression. We computed rate ratios (RR) using 285
weighted Cox proportional hazard models. Confi- 286
dence intervals were calculated using robust (or 287
sandwich) estimators of the variance. The use of ro- 288
bust estimators of the variance is recommended in 289
this case because the ordinary Wald confidence 290
intervals will not be guaranteed to provide at least 291
95% coverage probability [16]. 292

$$293 \quad SW = \frac{f(\text{BMI})}{f(\text{BMI}|\text{age, sex, physical activity, alcohol, smoking, fruit \& vegetable consumption, marital status})} \quad (1)$$

294 population, under the assumption of conditional
295 exchangeability [4, 6, 16]. For a formal definition of
296 causal effects and an explanation of the assumption

To adjust for selection bias, we computed another 293
set of stabilized weights, with the probability of being 294
non-censored given BMI in the numerator, and the 295

296 probability of being non-censored, given BMI and
 297 other variables associated with censoring status and
 298 the outcome (Equation 2; $C = 0$ for those non-
 299 censored) [4]. The numerator could be computed di-
 300 rectly from the data as the retention proportion by
 301 BMI categories. The denominator was computed
 302 with a logistic regression with non-censoring as the
 303 outcome and factors associated to censoring as
 304 independent variables.

$$SW' = \frac{f(C=0|BMI)}{f(C=0|BMI, \text{age, sex, physical activity, alcohol, smoking, fruit \& vegetable consumption, marital status})} \quad (2)$$

305 The final model used weights computed as the
 306 product of the stabilized weights for confounding
 307 adjustment times the stabilized weights for selection
 308 bias adjustment.

$$\text{Finalweights} = SW' \times SW \quad (3)$$

311 In the absence of unmeasured confounding,
 312 unmeasured informative censoring, and model mis-
 313 specification, the use of final weights (Equation 3)
 314 creates a pseudopopulation that corrects both selec-
 315 tion bias and confounding [19]. Analyses were per-
 316 formed with SAS version 9 (SAS Institute, Cary, NC,
 317 USA). Weighted Cox proportional hazard models
 318 were run using the WEIGHT statement, and the
 319 COVS options in the PROC PHREG statement.

321 Finally, to assess if selection bias depended on the
 322 main exposure of interest, we repeated the IPW anal-
 323 ysis to analyze the association between age at baseline
 324 (categorized as < 25, 25–34, 35–44, ≥45 years) and the
 325 risk of HTN. Models to compute weights included the
 326 same variables as in the BMI case.

327 Results

328 At the time of the analysis, 8647 of the 9907 eligible to
 329 have submitted their 2-year follow up questionnaire
 330 had done so (87%). Table 1 shows some characteris-
 331 tics of respondents and non-respondents to the follow-
 332 up questionnaire. In general, non-respondents had a
 333 higher proportion of individuals with missing values at
 334 the baseline questionnaire for all variables considered.
 335 Among non-respondents, age and fruit and vegetable
 336 consumption were lower, whereas the proportion of
 337 smokers, of non-married participants and of those
 338 without prior history of disease was higher.

339 In the univariate logistic regression analysis
 340 (Table 1) being younger, non-married, smoker, not
 341 having a driver's license and having suffered a pre-
 342 vious motor vehicle accident with hospitalization
 343 were associated with a higher probability of attrition.
 344 On the other hand, a previous history of cardiovas-
 345 cular disease or hypercholesterolemia at baseline
 346 increased the probability of responding the 2-year

follow-up questionnaire. Having a missing value in
 the baseline questionnaire in any item was associated
 with a higher probability of attrition.

When we adjusted simultaneously for all the
 variables shown in Table 1, a younger age, being
 non-married, smoker, obese, having a history of
 motor vehicle accident with hospitalization and an
 absence of previous history of cardiovascular disease
 or injury were associated with a higher probability

of attrition. The Hosmer-Lemeshow did not indicate
 a statistically significant lack of fit of our model.
 Only 3% of our sample would be misclassified in the
 contingency table, which is used to compute the
 Hosmer-Lemeshow test (comparing observed versus
 expected counts across deciles of the predicted
 probability).

To study the association between BMI and the risk
 of HTN, from the initial 9907 study participants,
 3321 participants were excluded due to either having
 prevalent HTN, cardiovascular disease, cancer or
 diabetes, their reporting of extreme caloric intakes or
 presenting missing values in some variables. That left
 6686 for the analysis, 5880 (88%) of whom had sub-
 mitted their first follow-up questionnaire. There were
 180 new cases of HTN during 13,526 person-years of
 follow-up. In Table 2, we show the rate ratios of
 HTN according to BMI categories. In the crude
 analysis, BMI was associated with an important in-
 crease in the risk of HTN. When we adjusted using
 IPW, the risk ratio estimates were attenuated but still
 showed an important direct association between BMI
 and the risk of HTN. Age and sex were the variables
 that accounted for most of the change in the risk
 ratios between crude and multivariate estimations.
 Additional adjustment for selection bias did not
 substantially change the risk ratio estimates. We re-
 peated the IPW analysis considering age as the main
 exposure, and we did not observe an important effect
 of censoring due to dropout in the estimates (data not
 shown).

Discussion

In this assessment of the first 2-years follow-up of
 the first 10,000 participants of the SUN cohort, we
 show that, as it happens with similar cohorts in the
 US, follow-up through mailed questionnaires of a
 highly educated cohort in Spain is feasible. Com-
 parable populations in other European countries
 could be a valuable setting to conduct studies with
 similar design. We also show that individuals who
 were lost to follow-up were different than those

Table 1. Baseline characteristics of the first 9907 participants in the SUN cohort by respondent status to the 2-year follow-up questionnaire, and odds ratios and 95% confidence intervals of loss to follow-up according to these characteristics

	Respondents (n = 8647)*	Non-respondents (n = 1260)*	Total (n = 9907)	Crude OR	Adjusted OR ^a
<i>Sociodemographic variables</i>					
<i>Age</i>					
< 25	1363 (83.1)	277 (16.9)	1640	1 (ref.)	1 (ref.)
25–34	2970 (86.1)	481 (13.9)	3451	0.8 (0.7–0.9)	0.8 (0.7–1.0)
35–44	2088 (92.1)	180 (7.9)	2268	0.4 (0.3–0.5)	0.5 (0.4–0.7)
≥45	2201 (89.5)	259 (10.5)	2460	0.6 (0.5–0.7)	0.5 (0.4–0.7)
Missing	25 (28.4)	63 (71.6)	88	12.4 (7.7–20.1)	N/A
<i>Gender</i>					
Female	5131 (87.1)	761 (12.9)	5892	1 (ref.)	1 (ref.)
Male	3439 (88.0)	469 (12.0)	3908	0.9 (0.8–1.0)	1.0 (0.9–1.2)
Missing	77 (72.0)	30 (28.0)	107	2.6 (1.7–4.0)	N/A
<i>Marital status</i>					
Married	4131 (90.0)	46 (10.0)	4592	1 (ref.)	1 (ref.)
Single	4060 (85.3)	701 (14.7)	4761	1.5 (1.4–1.8)	1.1 (1.0–1.3)
Widowed	98 (81.0)	23 (19.0)	121	2.1 (1.3–3.3)	2.4 (1.4–3.9)
Separated/divorced/others	197 (83.8)	38 (16.2)	235	1.7 (1.2–2.5)	1.7 (1.2–2.6)
Missing	161 (81.3)	37 (18.7)	198	2.1 (1.4–3.0)	N/A
<i>Lifestyle factors</i>					
<i>Physical activity^b</i>					
Sedentary	1863 (86.0)	303 (14.0)	2166	1 (ref.)	1 (ref.)
Less active	2268 (87.5)	324 (12.5)	2592	0.9 (0.7–1.0)	1.0 (0.8–1.2)
Moderately active	2379 (88.7)	303 (11.3)	2682	0.8 (0.7–0.9)	0.8 (0.7–1.0)
Very active	2062 (87.2)	303 (12.8)	2365	0.9 (0.8–1.1)	1.0 (0.8–1.2)
Missing	75 (73.5)	27 (26.5)	102	2.2 (1.4–3.5)	N/A
<i>Smoking</i>					
Never smoker	3824 (87.9)	528 (12.1)	4352	1 (ref.)	1 (ref.)
Past smoker	2226 (88.9)	279 (11.1)	2505	0.9 (0.8–1.1)	1.1 (0.9–1.3)
Current smoker	2226 (85.2)	386 (14.8)	2612	1.3 (1.1–1.4)	1.3 (1.1–1.5)
Missing	371 (84.7)	67 (15.3)	438	1.3 (1.0–1.7)	N/A
<i>Use of seatbelt</i>					
Always	7085 (87.3)	1028 (12.7)	8113	1 (ref.)	1 (ref.)
Not always	1240 (88.8)	157 (11.2)	1397	0.9 (0.7–1.0)	0.9 (0.8–1.1)
Almost never	168 (84.0)	32 (16.0)	200	1.3 (0.9–1.9)	1.1 (0.7–1.7)
Missing	154 (78.2)	43 (21.8)	197	1.9 (1.4–2.7)	N/A
<i>Use of alcohol and driving</i>					
Never	3748 (87.2)	550 (12.8)	4298	1 (ref.)	1 (ref.)
Almost never	1506 (88.4)	197 (11.6)	1703	0.9 (0.8–1.1)	0.9 (0.7–1.1)
Sometimes	2456 (89.0)	303 (11.0)	2759	0.8 (0.7–1.0)	0.8 (0.7–1.0)
Do not know how to drive	775 (82.2)	168 (17.8)	943	1.5 (1.2–1.8)	1.3 (1.1–1.6)
Missing	162 (79.4)	42 (20.6)	204	1.8 (1.2–2.5)	N/A
<i>Alcohol consumption (g/day)</i>					
Non-drinker	2011 (87.9)	278 (12.1)	2289	1 (ref.)	1 (ref.)
Drinker					
< 3 g/d	2476 (87.7)	346 (12.3)	2822	1.0 (0.9–1.2)	1.1 (0.9–1.3)
3–12 g/d	2689 (86.9)	406 (13.1)	3095	1.1 (0.9–1.3)	1.1 (0.9–1.4)
≥12 g/d	1397 (87.3)	203 (12.7)	1600	1.1 (0.9–1.3)	1.2 (1.0–1.5)
Missing	74 (73.3)	27 (26.7)	101	2.6 (1.7–4.2)	N/A
<i>Fruit and vegetable consumption</i>					
Below median	4252 (86.7)	652 (13.3)	4904	1 (ref.)	1 (ref.)
Above median	4321 (88.1)	581 (11.9)	4902	0.9 (0.8–1.0)	0.9 (0.8–1.0)
Missing	74 (73.3)	27 (26.7)	101	2.4 (1.5–3.7)	N/A
<i>Fat (% energy intake)</i>					
Below median	4292 (87.4)	617 (12.6)	4909	1 (ref.)	1 (ref.)
Above median	4273 (87.4)	614 (12.6)	4887	1.0 (0.9–1.1)	0.9 (0.8–1.0)
Missing	82 (73.9)	29 (26.1)	111	2.5 (1.6–3.8)	N/A
<i>Protein (% energy intake)</i>					
Below median	4243 (87.6)	598 (12.4)	4841	1 (ref.)	1 (ref.)

Table 1. Continued

	Respondents (n = 8647)*	Non-respondents (n = 1260)*	Total (n = 9907)	Crude OR	Adjusted OR ^a
Above median	4322 (87.2)	633 (12.8)	4955	1.0 (0.9-1.2)	1.1 (1.0-1.2)
Missing	82 (73.9)	29 (26.1)	111	2.5 (1.6-3.9)	N/A
<i>Clinical variables</i>					
BMI (kg/m²)					
< 25	5962 (87.4)	856 (12.6)	6818	1 (ref.)	1 (ref.)
25-30	1962 (87.8)	273 (12.2)	2235	1.0 (0.8-1.1)	1.1 (1.0-1.4)
> 30	363 (84.8)	65 (15.2)	428	1.2 (0.9-1.6)	1.5 (1.1-2.0)
Missing	360 (84.5)	66 (15.5)	426	1.3 (1.0-1.7)	N/A
Cardiovascular disease					
No	8264 (87.1)	1225 (12.9)	9489	1 (ref.)	1 (ref.)
Yes	383 (91.6)	35 (8.4)	418	0.6 (0.4-0.9)	0.6 (0.4-0.9)
Hypercholesterolemia					
No	7323 (86.9)	1103 (13.1)	8426	1 (ref.)	1 (ref.)
Yes	1324 (89.4)	157 (10.6)	1481	0.8 (0.7-0.9)	1.0 (0.8-1.2)
Diabetes					
No	8505 (87.2)	1246 (12.8)	9751	1 (ref.)	1 (ref.)
Yes	142 (91.0)	14 (9.0)	156	0.7 (0.4-1.2)	0.7 (0.4-1.3)
Hypertension					
No	7842 (87.3)	1141 (12.7)	8983	1 (ref.)	1 (ref.)
Yes	805 (87.1)	119 (12.9)	924	1.0 (0.8-1.2)	1.2 (1.0-1.5)
Motor vehicle accident with hospitalization					
No	8237 (87.5)	1180 (12.5)	9417	1 (ref.)	1 (ref.)
Yes	410 (83.7)	80 (16.3)	490	1.4 (1.1-1.7)	1.4 (1.1-1.9)

N/A: not applicable. *Values are numbers (percentage, calculated from total column).

^aAdjusting simultaneously for all variables in the table but excluding 925 individuals with missing values in any variable in the table.

^bPhysical activity was categorized as sedentary if the individual reported no leisure-time physical activity, and less, moderate or very active using tertiles of MET-hours/week.

Table 2. Rate ratios (95% confidence intervals) of incident hypertension in the SUN cohort according to body mass index categories

Body mass index (kg/m ²)	< 25	25-30	> 30
Cases of HTN	76	84	20
Person-years	10219.0	2880.7	426.2
Crude RR	1 (ref.)	4.0 (2.9-5.4)	6.4 (3.9-10.5)
Adjusting for censoring using IPW ^a	1 (ref.)	4.0 (2.9-5.5)	6.7 (4.0-11.2)
Adjusting for confounding using IPW ^b	1 (ref.)	2.1 (1.5-3.0)	2.4 (1.1-5.3)
Adjusting for both confounding and selection bias using IPW ^c	1 (ref.)	2.1 (1.4-3.0)	2.4 (1.1-5.3)

HTN: hypertension; IPW: inverse probability weighting; RR: rate ratio.

^aAdjusted for selection bias. Weights computed according to formula 2 in the text.

^bAdjusted for age, sex, physical activity, alcohol intake, smoking, fruit and vegetable consumption, and marital status. Weights computed according to formula 1 in the text.

^cAdjusted for age, sex, physical activity, alcohol intake, smoking, fruit and vegetable consumption, and marital status, and for selection bias. Weights computed according to formula 3 in the text.

397 retained in the study. In fact, non-respondents were
398 more likely to be younger, not married, obese or
399 overweight, smokers and to belong to the group not
400 having a driving license. In contrast, a previous
401 history of cardiovascular disease increased the
402 probability of being a respondent. Finally, we show
403 that censoring due to attrition did not introduce an
404 important selection bias for the association between
405 BMI and HTN when we corrected this bias using
406 IPW.

Several reasons could explain why different vari-
ables predicted lost to follow-up in the study. Age is
a factor strongly correlated with censoring in lon-
gitudinal studies. On the one hand, younger people
tend to move more frequently, and that results in a
greater difficulty to locate them [9, 20]. On the other
hand, older age has been associated consistently
with higher dropout rates in cohort studies con-
ducted in the elderly [3, 21]. Our study population is
relatively young (mean age 36 at baseline), and

probably very mobile, causing the tracking and retention of participants more difficult through mailed questionnaires. In the Cornellá Health Interview Survey Follow-up Study (another Spanish cohort), being single was an important correlate of migration out of the study area [20], which is consistent with our results. In our study, married people had the higher retention rate.

Smoking and obesity were more frequent among participants lost to follow-up than among retained individuals. Other studies have found similar results regarding these two variables [22, 23]. Then, it can be said that health-conscious individuals have lower dropout probabilities in epidemiologic studies. The other factor associated with a higher retention rate was a previous history of cardiovascular disease. Probably, the experience of having been diagnosed of a cardiovascular disease causes these individuals to be more health conscious than the rest of the cohort, and more willing to keep participating in the cohort.

Finally, a higher educational level is one of the strongest and more consistent predictors for retention in cohort studies [9, 20, 22, 23]. In fact, this is one of the reasons to explain the decision to restrict the SUN study to university graduates. Interestingly, the retention in the cohort was similar for those with a postgraduate degree (master, PhD) than for those who had attained only a college degree.

Different methods to correct for selection bias due to differential retention have been proposed, mainly based on imputation methods for lost to follow-up individuals [24]. In this study, we have used IPW as a method to adjust both for confounding and selection bias. The use of IPW allows the estimation of causal effect under the assumption of unmeasured confounding and unmeasured informative censoring. To compute weights, we have taken into consideration the main risk factors for HTN that could act as confounding factors [25], and we can therefore be confident on the appropriateness of both assumptions.

Adjustment for confounding using IPW is not enough to remove selection bias caused by dropouts (although a traditional multivariate analysis would be sufficient in some cases). This is because in the pseudopopulation created applying weights to adjust for confounding (formula 1), there would not be any association between the potential confounders included in the IPW calculations and the exposure of interest but the possible association between those variables and censoring status, and the selection bias it causes, would remain. Therefore, additional IPWs are needed to take into account both censoring and confounding [4].

Finally, although we found no apparent selection bias for the association between BMI and HTN, or between age and HTN, we cannot be sure that these results can be generalized to other outcomes and exposures. In fact, some factors affecting

probability of retention could be independent of the incidence of HTN (and then would not create selection bias), but they might be associated with another outcome. In that case, weight calculations for IPW adjustment should include that factor. In general, decisions to include or not a given factor in any weight calculation should be based in a priori subject matter knowledge [26]. Overall, study of selection bias due to censoring should be performed for each particular instance of an exposure-outcome relationship.

In conclusion, we show that the follow-up through mailed questionnaires of a geographically disperse cohort in Spain is possible. There are differences between retained participants and individuals lost to follow-up, but this does not necessarily have an important impact in the rate ratio estimates.

Acknowledgements

The authors are indebted to Dr Miguel Hernán for comments on this manuscript and Ms Carmen de la Fuente for technical assistance. Also, they thank the continuous collaboration of the participants in the SUN Study. The SUN Study has received funding from the Spanish Ministry of Health (Grants PI040233, G03/140, and PI030678), the Navarra Regional Government (43/2002 and 41/2005) and the University of Navarra. Dr Alonso was supported partially by a Fulbright fellowship and a MMA Foundation grant.

References

- Delgado-Rodríguez M, Llorca J. Bias. *J Epidemiol Community Health* 2004; 58: 635-641.
- Greenland S. Response and follow-up bias in cohort studies. *Am J Epidemiol* 1977; 106: 184-187.
- Chatfield MD, Brayne CE, Matthews FE. A systematic literature review of attrition between waves in longitudinal studies in the elderly shows a consistent pattern of dropout between differing studies. *J Clin Epidemiol* 2005; 58: 13-19.
- Hernán MA, Hernández-Díaz S, Robins JM. A structural approach to selection bias. *Epidemiology* 2004; 15: 615-625.
- Greenland S. Basic methods for sensitivity analysis and external adjustment. In: Rothman KJ and Greenland S (eds.) *Modern Epidemiology*. Philadelphia: Lippincott Williams and Wilkins, 1998, p. 343-357.
- Hernán MA, Robins JM. Estimating causal effects from epidemiologic data. *J Epidemiol Community Health* 2006; (in press).
- Colditz GA, Manson JE, Hankinson SE. The Nurses' Health Study: 20-year contributing to the understanding of health among women. *J Womens Health* 1997; 6: 49-62.



531	8. Koh-Banerjee P, Wang Y, Hu FB, Spiegelman D, Willett WC, Rimm EB. Changes in body weight and body fat distribution as risk factors for clinical diabetes in US men. <i>Am J Epidemiol</i> 2004; 159: 1150-1159.	571
532		572
533		573
534		574
535	9. Russell C, Palmer JR, Adams-Campbell LL, Rosenberg L. Follow-up of a large cohort of black women. <i>Am J Epidemiol</i> 2001; 154: 845-853.	575
536		576
537		577
538	10. Hunt JR, White E. Retaining and tracking cohort study members. <i>Epidemiol Rev</i> 1998; 20: 57-70.	578
539		579
540	11. Martínez-González MA, Sánchez-Villegas A, de Irala-Estévez J, Martí A, Martínez JA. Mediterranean diet and stroke: Objectives and design of the SUN Project. <i>Nutr Neurosci</i> 2002; 5: 65-73.	580
541		581
542		582
543		583
544	12. Martín-Moreno JM, Boyle P, Gorgojo L, et al. Development and validation of a food frequency questionnaire in Spain. <i>Int J Epidemiol</i> 1993; 22: 512-519.	584
545		585
546		586
547	13. Ainsworth BE, Haskell WL, Whitt MC, et al. Compendium of physical activities: An update of activity codes and MET intensities. <i>Med Sci Sports Exerc</i> 2000; 32: S498-S504.	587
548		588
549		589
550		590
551	14. Bes-Rastrollo M, Pérez Valdivieso JR, Sánchez-Villegas A, Alonso A, Martínez-González MA. Validación del peso e índice de masa corporal auto-declarados de los participantes de una cohorte de graduados universitarios. <i>Rev Esp Obes</i> 2005; (in press).	591
552		592
553		593
554		594
555		595
556	15. Alonso A, Beunza JJ, Delgado-Rodríguez M, Martínez-González MA. Validation of self reported diagnosis of hypertension in a cohort of university graduates in Spain. <i>BMC Public Health</i> 2005; 5: 94.	596
557		597
558		598
559		599
560	16. Robins JM, Hernán MA, Brumback B. Marginal structural models and causal inference in epidemiology. <i>Epidemiology</i> 2000; 11: 550-560.	600
561		601
562		602
563	17. World Health Organization. <i>Obesity: Preventing and Managing the Global Epidemic</i> . Geneva: World Health Organization, 2000.	
564		
565		
566	18. Hernán MA. A definition of causal effect for epidemiologic research. <i>J Epidemiol Community Health</i> 2004; 58: 265-271.	
567		
568		
569	19. Cole SR, Hernán MA, Robins JM, et al. Effect of highly active antiretroviral therapy on time to acquired immunodeficiency syndrome or death using marginal structural models. <i>Am J Epidemiol</i> 2003; 158: 687-694.	
570		
	20. García M, Fernández E, Schiaffino A, Borrell C, Martí M, Borrás JM. Attrition in a population-based cohort eight years after baseline interview: The Cornella Health Interview Survey Follow-up (CHIS.FU) Study. <i>Ann Epidemiol</i> 2004; 15: 98-104.	
	21. Zunzunegui MV, Beland F, Gutierrez Cuadra P. Loss to follow-up in a longitudinal study on aging in Spain. <i>J Clin Epidemiol</i> 2001; 54: 501-510.	
	22. Matthews FE, Chatfield MD, Freeman C, McCracken C, Brayne CE, MRC CFAS. Attrition and bias in the MRC cognitive function and ageing study: An epidemiological investigation. <i>BMC Public Health</i> 2004; 4: 12-21.	
	23. Pirzada A, Yan LL, Garside DB, Schiffer L, Dyer AR, Daviglius ML. Response rates to a questionnaire 26 years after baseline examination with minimal interim participant contact and baseline differences between respondents and nonrespondents. <i>Am J Epidemiol</i> 2004; 159: 94-101.	
	24. Twisk J, de Vente W. Attrition in longitudinal studies. How to deal with missing data. <i>J Clin Epidemiol</i> 2002; 55: 329-337.	
	25. Whelton PK, He J, Appel LJ, et al. Primary prevention of hypertension: Clinical and public health advisory from the National High Blood Pressure Education Program. <i>JAMA</i> 2002; 288: 1882-1888.	
	26. Hernán MA, Hernández-Díaz S, Werler MM, Mitchell AA. Causal knowledge as a prerequisite for confounding evaluation: An application to birth defects epidemiology. <i>Am J Epidemiol</i> 2002; 155: 176-184.	
	<i>Address for correspondence:</i> Alvaro Alonso, Department of Epidemiology, Harvard School of Public Health, 677 Huntington Avenue, Boston, MA 02115, USA Phone: +1-617-432-2838; Fax: +1-617-566-7805 E-mail: aalogut@alumni.unav.es	
		603
		604
		605
		606
		607
		608